

Basic Phylogenetic Combinatorics

Phylogenetic combinatorics is a branch of discrete applied mathematics concerned with the combinatorial description and analysis of phylogenetic trees and related mathematical structures such as phylogenetic networks and tight spans. Based on a natural conceptual framework, the book focuses on the interrelationship between the principal options for encoding phylogenetic trees: split systems, quartet systems, and metrics. Such encodings provide useful options for analyzing and dealing with phylogenetic trees and networks, and are at the basis of much of phylogenetic data processing. The book highlights how each one provides a unique perspective for viewing and perceiving the combinatorial structure of a phylogenetic tree and is, simultaneously, a rich source for combinatorial analysis and theory building. It is dedicated to Manfred Eigen who inspired many of the results presented in this book.

Graduate students and researchers in mathematics and computer science will enjoy exploring this fascinating new area, and learn how mathematics may be used to help solve topical problems arising in evolutionary biology.

ANDREAS DRESS works currently as a scientific advisor at infinity³ GmbH, Bielefeld, Germany.

KATHARINA T. HUBER is a Lecturer in the School of Computing Sciences at the University of East Anglia, UK.

JACOBUS KOOLEN is an Associate Professor in the Department of Mathematics at Pohang University of Science and Technology (POSTECH), South Korea.

VINCENT MOULTON is a Professor in the School of Computing Sciences at the University of East Anglia, UK.

ANDREAS SPILLNER is an Assistant Professor in the Department of Mathematics and Computer Science at the University of Greifswald, Germany.

Basic Phylogenetic Combinatorics

ANDREAS DRESS (德乐思)

infinity³ GmbH, Bielefeld, Germany

and

*CAS-MPG Partner Institute for Computational
Biology, Shanghai Institutes for Biological Sciences*

KATHARINA T. HUBER

University of East Anglia

JACOBUS KOOLEN

*Pohang University of Science and Technology (POSTECH),
South Korea*

VINCENT MOULTON

University of East Anglia

ANDREAS SPILLNER

University of Greifswald



CAMBRIDGE
UNIVERSITY PRESS

Cambridge University Press & Assessment

978-0-521-76832-0 — Basic Phylogenetic Combinatorics

Andreas Dress , Katharina T. Huber , Jacobus Koolen , Vincent Moulton , Andreas Spillner
Frontmatter

[More Information](#)



CAMBRIDGE
UNIVERSITY PRESS

Shaftesbury Road, Cambridge CB2 8EA, United Kingdom

One Liberty Plaza, 20th Floor, New York, NY 10006, USA

477 Williamstown Road, Port Melbourne, VIC 3207, Australia

314–321, 3rd Floor, Plot 3, Splendor Forum, Jasola District Centre, New Delhi – 110025, India

103 Penang Road, #05–06/07, Visioncrest Commercial, Singapore 238467

Cambridge University Press is part of Cambridge University Press & Assessment,
a department of the University of Cambridge.

We share the University's mission to contribute to society through the pursuit of
education, learning and research at the highest international levels of excellence.

www.cambridge.org

Information on this title: www.cambridge.org/9780521768320

© A. Dress, K. T. Huber, J. Koolen, V. Moulton and A. Spillner 2012

This publication is in copyright. Subject to statutory exception and to the provisions
of relevant collective licensing agreements, no reproduction of any part may take
place without the written permission of Cambridge University Press & Assessment.

First published 2012

A catalogue record for this publication is available from the British Library

Library of Congress Cataloging-in-Publication data

Basic phylogenetic combinatorics / Andreas Dress . . . [et al.].

p. cm.

ISBN 978-0-521-76832-0 (Hardback)

1. Branching processes. 2. Combinatorial analysis. I. Dress, Andreas.

QA274.76.B37 2011

511'.6–dc23

2011043264

ISBN 978-0-521-76832-0 Hardback

Cambridge University Press & Assessment has no responsibility for the persistence
or accuracy of URLs for external or third-party internet websites referred to in this
publication and does not guarantee that any content on such websites is, or will
remain, accurate or appropriate.

We dedicate this book to Manfred Eigen whose questions concerning the evolution of RNA and DNA sequences inspired many of the early results in phylogenetic combinatorics that ultimately led to the work presented in this book.

Contents

	<i>Preface</i>	<i>page ix</i>
1	Preliminaries	1
	1.1 Sets, set systems, and partially ordered sets	1
	1.2 Graphs	4
	1.3 Metric spaces	13
	1.4 Computational complexity	19
2	Encoding X-trees	21
	2.1 X -trees	21
	2.2 Encoding X -trees with splits	23
	2.3 Encoding X -trees with metrics	26
	2.4 Encoding X -trees with quartets	27
3	Consistency of X-tree encodings	31
	3.1 The 4-point condition	31
	3.2 Compatibility	38
	3.3 Quartet systems	42
4	From split systems to networks	50
	4.1 The Buneman graph	51
	4.2 The Buneman graph of a compatible split system	59
	4.3 Median networks	63
	4.4 Split networks	65
	4.5 Split graphs and metrics: The theory of X -nets	72
5	From metrics to networks: The tight span	75
	5.1 The tight span	75
	5.2 A canonical contraction from $P(D)$ onto $T(D)$	82
	5.3 The tight span of a finite metric space	87
	5.4 Networks from tight spans	93

5.5	Network realizations of metrics	97
5.6	Optimal and hereditarily optimal realizations	100
6	From quartet and tree systems to trees	104
6.1	On quartet systems	105
6.2	On set and tree systems	113
6.3	Constructing trees from quartet, tree, and set systems	118
6.4	Slim tree systems	121
6.5	Definitive set systems	128
7	From metrics to split systems and back	137
7.1	Buneman splits	137
7.2	Weakly compatible split systems	146
7.3	From weighted split systems to bivariate maps	161
7.4	The Buneman complex and the tight span	167
8	Maps to and from quartet systems	171
8.1	A Galois connection between split and quartet systems	171
8.2	A map from quartets to metrics	177
8.3	Transitive quartet systems	180
9	Rooted trees and the Farris transform	195
9.1	Rooted X -trees, clusters, and triplets	198
9.2	Dated rooted X -trees and hierarchical dissimilarities	202
9.3	Affine versus projective clustering and the combinatorial Farris transform	205
9.4	Hierarchical dissimilarities, hyperbolic maps, and their Farris transform	209
9.5	Hierarchical dissimilarities, generalized metrics, and the tight-span construction	214
9.6	Algorithmic issues	218
10	On measuring and removing inconsistencies	222
10.1	k -compatibility	222
10.2	Δ -hierarchical approximations	230
10.3	Quartet-Joining and QNet	236
	<i>Commonly used symbols</i>	242
	<i>Bibliography</i>	253
	<i>Index</i>	261

Preface

More than one and a half centuries have passed since Charles Darwin presented his theory on the origin of species asserting that all organisms are related to each other by common descent via a “tree of life”. Since then, biologists have been able to piece together a great deal of information concerning this tree — relying in particular in more recent times on the advent of ever cheaper and faster DNA sequencing technologies. Even so, there remain many fascinating open problems concerning the tree of life and the evolutionary processes underlying it, problems that often require sophisticated techniques from areas such as mathematics, computer science, and statistics.

Phylogenetic combinatorics can be regarded as a branch of discrete applied mathematics concerned with the combinatorial description and analysis of *phylogenetic* or *evolutionary trees* and related mathematical structures such as phylogenetic networks, complexes, and tight spans. In this book, we present a *systematic* approach to phylogenetic combinatorics based on a natural conceptual framework that, simultaneously, allows and forces us to encompass many classical as well as a good number of new pertinent results.

More specifically, this book concentrates on the interrelationship between the three principal ways commonly used for **encoding** phylogenetic trees: *Split systems*, *metrics*, and *quartet systems* (see Figure 1). Informally, for X some finite set, a split system over X is a collection of bipartitions of X , a quartet system is a collection of two-versus-two bipartitions of subsets of X of size four, and a metric is a bivariate function assigning a “distance” to any pair of elements in X .

Such encodings provide useful options for analyzing and manipulating phylogenetic trees with leaves labeled by X , and are at the basis of much of phylogenetic data processing. Indeed, they arise naturally from the various types of data from which phylogenetic trees are typically (re-)constructed: Comparative sequence analysis of genes or genomes may lead to metrics, character tables as

Preface

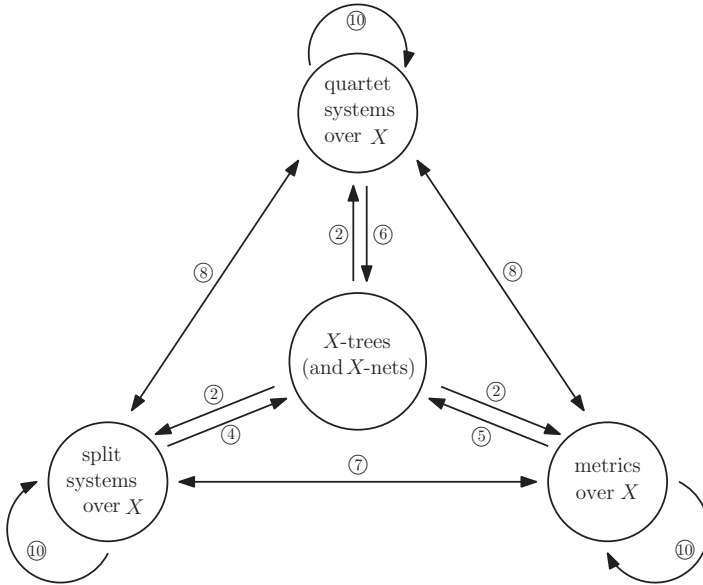


Figure 1 In this figure, we indicate the manifold relationships between various combinatorial objects relevant in phylogenetic analysis that will be studied in this book.

well as single nucleotide polymorphisms give rise to splits (and metrics), and quartet systems arise from restricting attention in phylogenetic data analysis to just four taxa at a time to avoid reconstruction algorithms becoming overwhelmed by the sheer number of taxa that need to be treated simultaneously. All three types of encodings require a solid theoretical foundation and provide, at the same time, a rich source for combinatorial analysis and theory building.

This book aims to highlight how each of the three types of encodings provides a unique perspective for viewing the combinatorial structure of a phylogenetic tree, for assessing the suitability of given data for tree reconstruction, and, if suitable, for recovering such trees from such data. And it will, of course, also discuss how split systems, metrics, and quartet systems are related to one another.

Here is an outline of the contents: After presenting some basic definitions and concepts that will be used throughout the book in Chapter 1, we introduce the formal concept of a phylogenetic tree or, a bit more generally, an X -tree in Chapter 2. We then define split systems, metrics, and quartet systems, and show that X -trees may indeed be uniquely encoded in terms of such combinatorial objects. In Chapter 3, we then proceed to identify which split systems, metrics,

or quartet systems are induced by — and thus encode — an X -tree: That is, we characterize the split systems, metrics, or quartet systems in the “image” of the “maps” labeled ② in Figure 1 in terms of some simple, yet instructive and enlightening conditions.

In Chapters 4, 5, and 6, we move on to the problem of deciding how to *decode* a given tree-encoding using appropriate constructions corresponding, respectively, to the maps labeled ④, ⑤, ⑥ in Figure 1. In other words, given a split system, a metric, or a quartet system in the image of the maps labeled ②, we consider how to find that (unique!) X -tree encoded by them. We will also explain how, when applied to data sets that do not encode a tree, these constructions can produce *networks* (as opposed to trees) and discuss some pertinent consequences.

In Chapters 7 and 8, we investigate the *recoding* problem: How can we compute the split system, metric, or quartet system encoding a tree from its other two encodings, i.e., how can we define maps (as indicated by the arrows labeled ⑦ and ⑧, in either direction) so that the resulting triangular subdiagrams in Figure 1 are commutative? Generally, the pertinent constructions lead to pairs of maps between any two of the three classes of objects that appear to be of some independent interest in themselves.

“Rooting” an X -tree is rather important for a realistic interpretation of X -trees in terms of evolutionary history. Correspondingly, we consider in Chapter 9 how the previously mentioned maps and constructions can be modified so as to apply to *rooted X-trees*, giving rise to *cluster systems* (rather than split systems), *triplet systems* (rather than quartet systems), and *hierarchical dissimilarities* and *ultrametrics* (rather than metrics). Mathematically speaking, this can be regarded as taking an *affine* (more concrete) versus a *projective* (more elegant) approach to working with phylogenetic trees.

In the final chapter, Chapter 10, we address the problem of how to measure and remove “inconsistencies” in split systems, metrics, and quartet systems. In other words, given one such structure that does not encode an X -tree, we explore how we may find another one in its “neighborhood” that does. As we shall see, this not only yields some interesting mathematical results, but also new ways to analyze and understand phylogenetic data.

A major feature of this book is that full proofs are provided for all of the fundamental results, thus giving the motivated reader a chance to get to the forefront of the field of phylogenetic combinatorics without having to spend too much time seeking references (see also Figure 2 for a *Leitfaden*, i.e., an overview of chapter dependencies). It also includes various new results and proofs that have not been published previously, and it attempts to introduce most topics in an elementary way. Overall, we hope that the reader will be

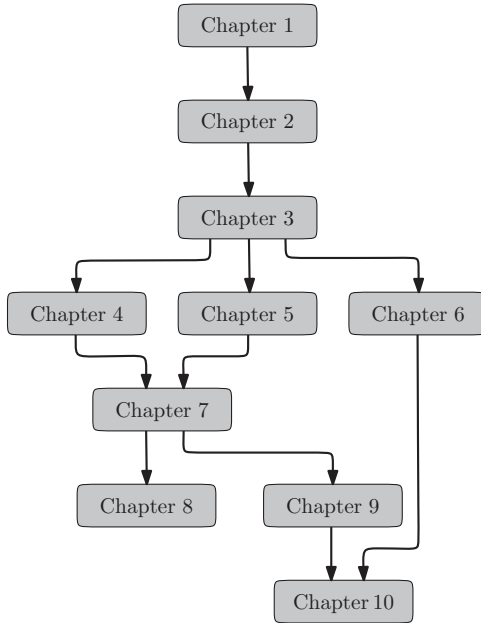
Preface

Figure 2 A diagram depicting the main dependencies between chapters in this book.

motivated by the book to explore this interesting area of mathematics whilst, at the same time, having the opportunity of seeing how mathematics may be used to help with solving topical problems that arise outside mathematics.

Finally, we would like to thank the Isaac Newton Institute for Mathematical Sciences, Cambridge, for giving us the opportunity to jointly draft and work on this book there, and also the UK Engineering and Physical Sciences Research Council and Royal Society, the Basic Science Research Program through the National Research Foundation of Korea (NRF) (grant number 2010-0008138), the DFG and the Max Planck Society, Germany, and the Chinese Academy for Sciences for financial support. We also thank our friends and colleagues and, in particular, David Bryant, Stefan Grünwald, Mike Hendy, Daniel Huson, Saitou Naruya, David Penny, LI Qiang, Charles Semple, Mike Steel, and WU Yaokun for many stimulating, critical, as well as encouraging discussions and comments. In addition, we thank students and colleagues at Bangalore, Bandar Lampung, Christchurch, Manila, New York, Paris, Pohang, and Shanghai, for their helpful feedback in courses where early versions of this text were presented. And last but not least, we all thank our families and, especially, Christiana, Keiko, Therese, Eugen, Jacky, and Robin, for their patience and support.